

Analyse topologique de données et estimation de support

EDDIE AAMARI

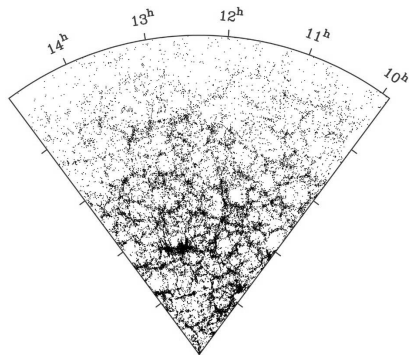
LABORATOIRE DE PROBABILITÉ, STATISTIQUES ET MODÉLISATION
CNRS, UNIVERSITÉ PARIS CITÉ, SORBONNE UNIVERSITÉ

CONGRÈS FRANÇAIS DE MÉCANIQUE

—
NANTES

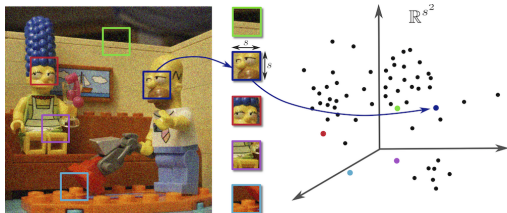
30 AOÛT 2022

Data with a Global Geometric Structure



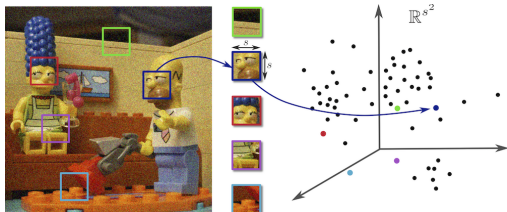
Large Scale Galaxy Structures: one point represents a galaxy in \mathbb{R}^3
[2dF Galaxy Redshift Survey]

Data with a Global Geometric Structure

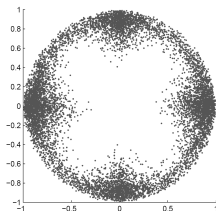


Extracting all the $s \times s$ patches of an image with $m \times n$ pixels.
[Houdard – 2018]

Data with a Global Geometric Structure

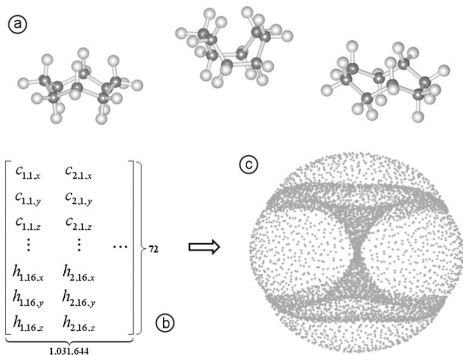


Extracting all the $s \times s$ patches of an image with $m \times n$ pixels.
[Houdard – 2018]



For $s = 7$, one image $M \in \mathbb{R}^{n \times m}$ yields $\asymp mn$ points in $\mathbb{R}^{7 \times 7} = \mathbb{R}^{49}$
[Xia – 2016]

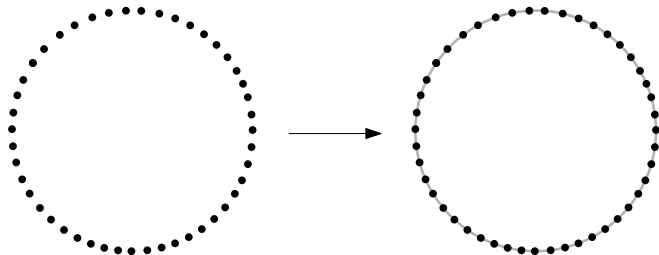
Data with a Global Geometric Structure



Cyclo-octane (C_8H_{16}) conformations
[Martin *et al.* – 2010]

One conformation is described with a point in $(\mathbb{R}^3)^{8+16} = \mathbb{R}^{72}$.

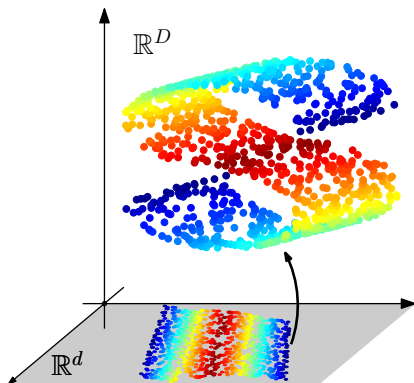
Uncover Data Structure



Input: a set $\mathbb{X}_n = \{X_1, \dots, X_n\}$ of observations.

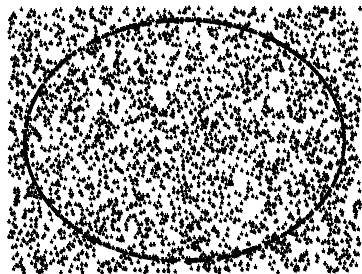
Goal: Understand the underlying structure of the data,
for interpretation or summary.

Challenge 1: Dimension



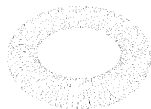
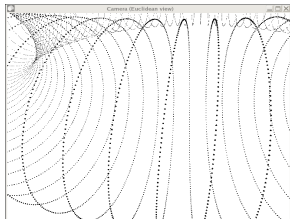
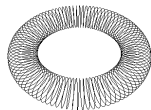
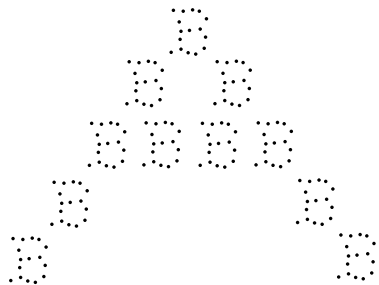
What dimension is this S-shape?

Challenge 2: Noise



Are my data corrupted?

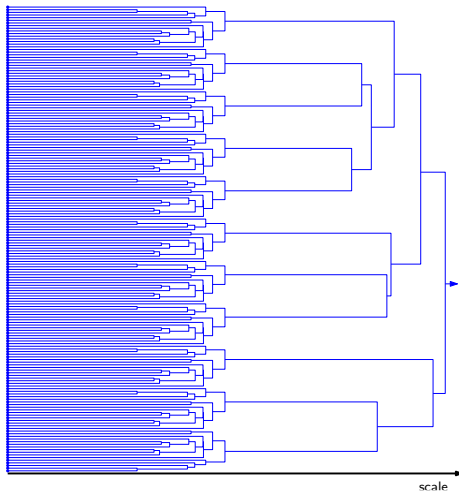
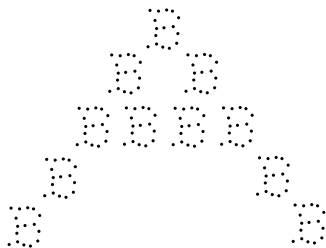
Challenge 3: Scale



Zoom in or zoom out?

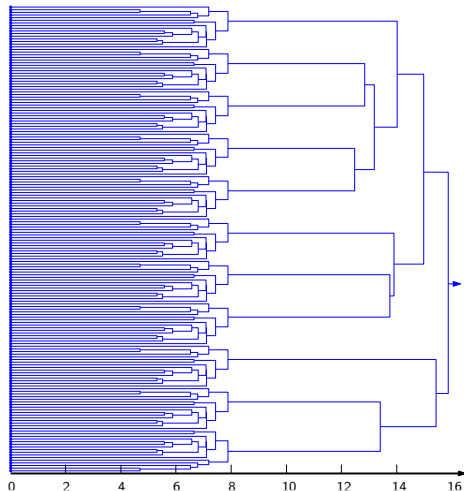
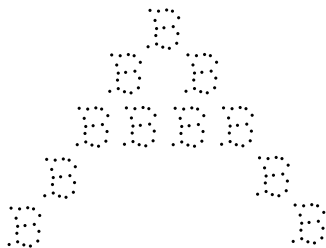
Dendrogram, Persistence Barcodes and Generalization

$$d_{\mathcal{P}} : \mathbb{R}^2 \rightarrow \mathbb{R}$$
$$x \mapsto \min_{p \in \mathcal{P}} \|x - p\|$$



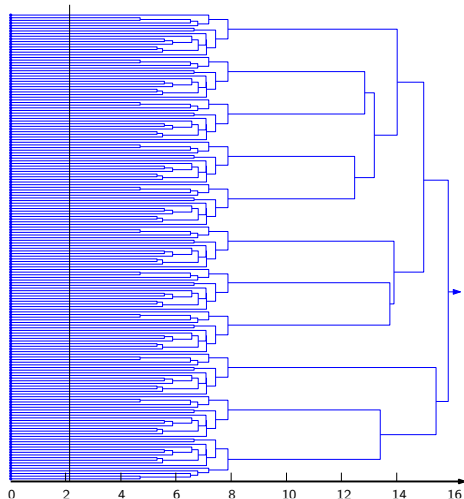
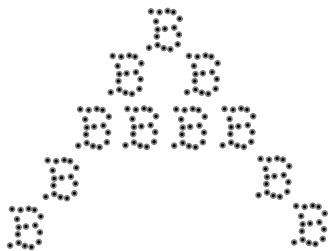
Dendrogram, Persistence Barcodes and Generalization

$$d_{\mathcal{P}} : \mathbb{R}^2 \rightarrow \mathbb{R}$$
$$x \mapsto \min_{p \in \mathcal{P}} \|x - p\|$$



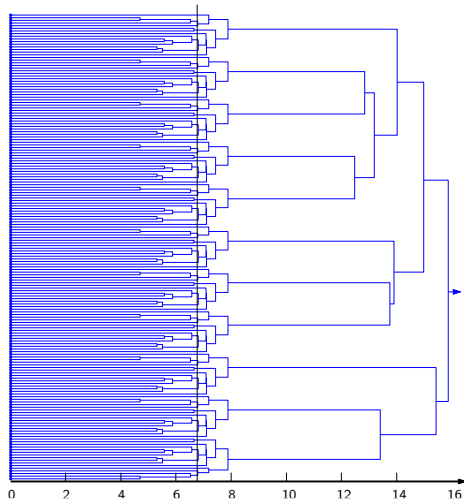
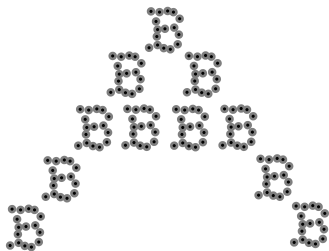
Dendrogram, Persistence Barcodes and Generalization

$$d_{\mathcal{P}} : \mathbb{R}^2 \rightarrow \mathbb{R}$$
$$x \mapsto \min_{p \in \mathcal{P}} \|x - p\|$$



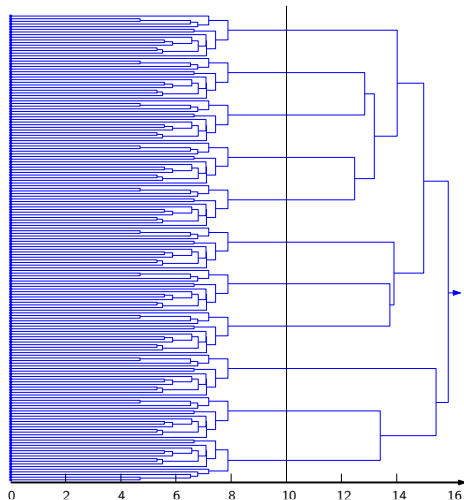
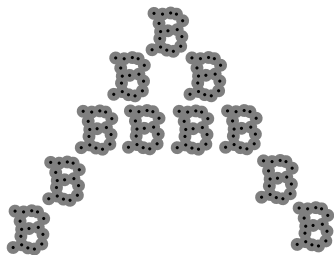
Dendrogram, Persistence Barcodes and Generalization

$$d_{\mathcal{P}} : \mathbb{R}^2 \rightarrow \mathbb{R}$$
$$x \mapsto \min_{p \in \mathcal{P}} \|x - p\|$$



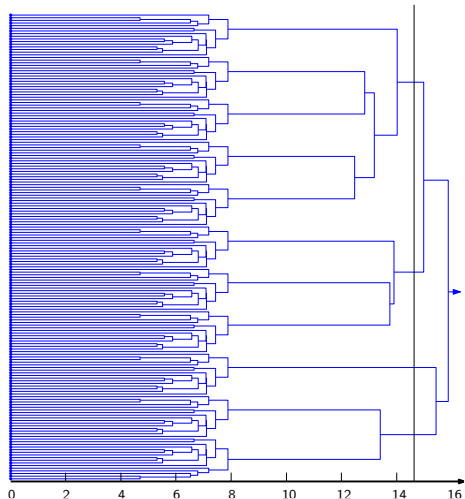
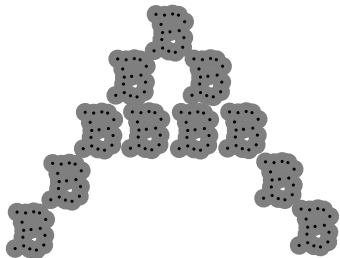
Dendrogram, Persistence Barcodes and Generalization

$$d_{\mathcal{P}} : \mathbb{R}^2 \rightarrow \mathbb{R}$$
$$x \mapsto \min_{p \in \mathcal{P}} \|x - p\|$$



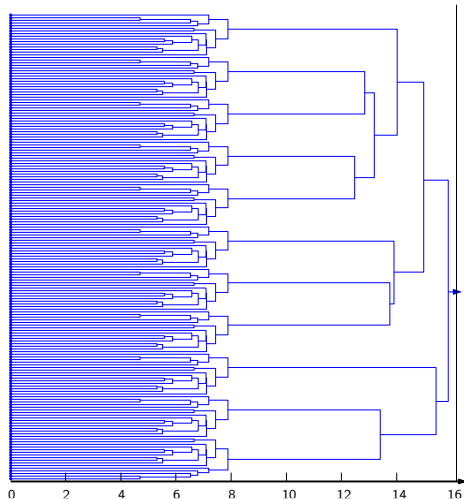
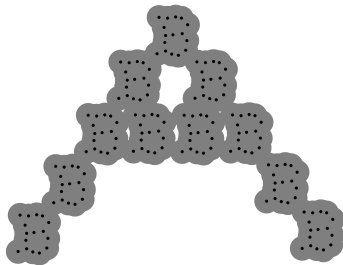
Dendrogram, Persistence Barcodes and Generalization

$$d_{\mathcal{P}} : \mathbb{R}^2 \rightarrow \mathbb{R}$$
$$x \mapsto \min_{p \in \mathcal{P}} \|x - p\|$$



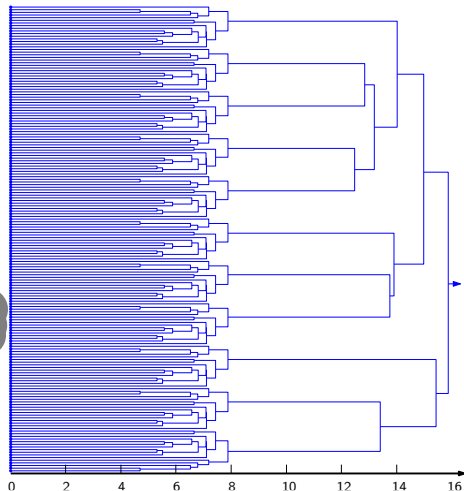
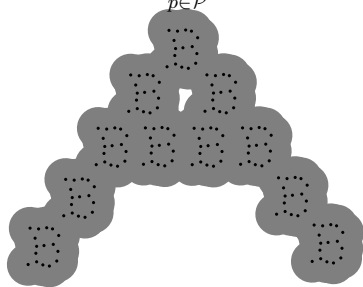
Dendrogram, Persistence Barcodes and Generalization

$$d_{\mathcal{P}} : \mathbb{R}^2 \rightarrow \mathbb{R}$$
$$x \mapsto \min_{p \in \mathcal{P}} \|x - p\|$$



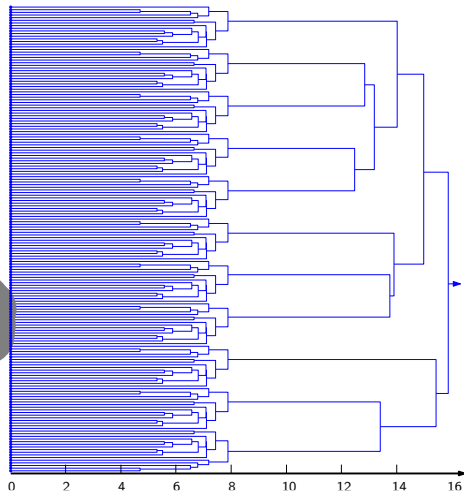
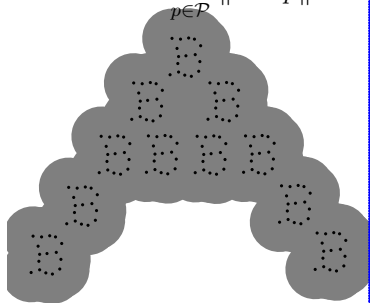
Dendrogram, Persistence Barcodes and Generalization

$$d_{\mathcal{P}} : \mathbb{R}^2 \rightarrow \mathbb{R}$$
$$x \mapsto \min_{p \in \mathcal{P}} \|x - p\|$$



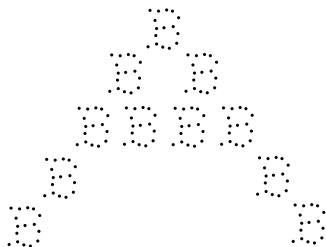
Dendrogram, Persistence Barcodes and Generalization

$$d_{\mathcal{P}} : \mathbb{R}^2 \rightarrow \mathbb{R}$$
$$x \mapsto \min_{p \in \mathcal{P}} \|x - p\|$$



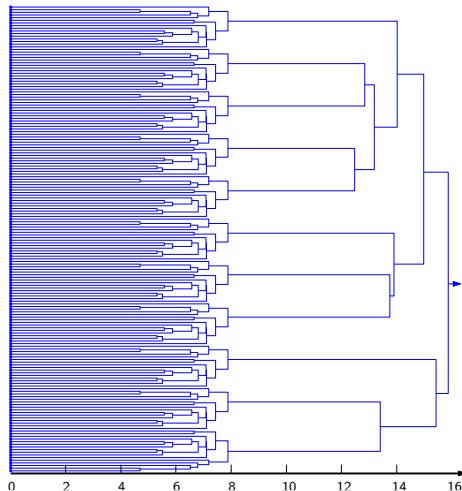
Dendrogram, Persistence Barcodes and Generalization

$$d_{\mathcal{P}} : \mathbb{R}^2 \rightarrow \mathbb{R}$$
$$x \mapsto \min_{p \in \mathcal{P}} \|x - p\|$$



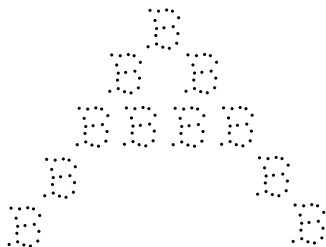
Dendrogram is:

- informative
- unstable

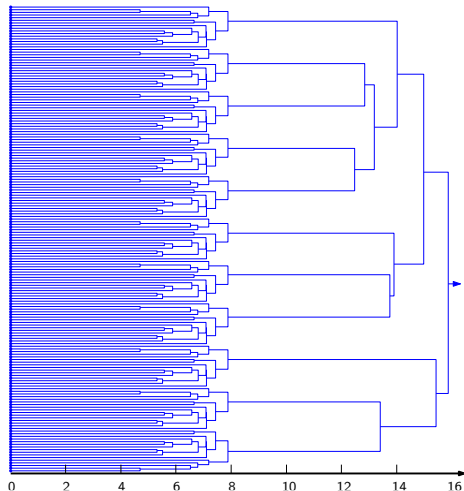


Dendrogram, Persistence Barcodes and Generalization

$$d_{\mathcal{P}} : \mathbb{R}^2 \rightarrow \mathbb{R}$$
$$x \mapsto \min_{p \in \mathcal{P}} \|x - p\|$$

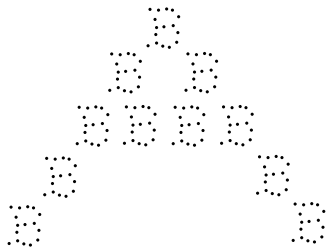


Dendrogram \rightarrow barcode

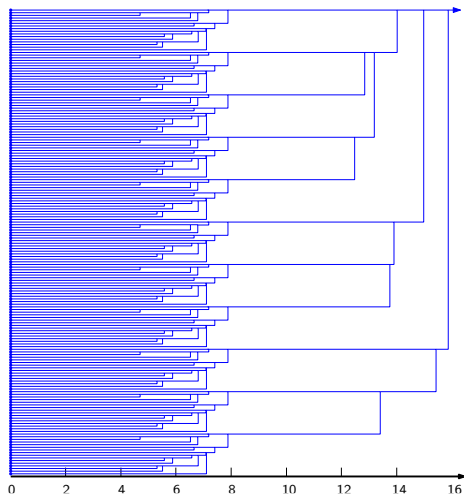


Dendrogram, Persistence Barcodes and Generalization

$$d_{\mathcal{P}} : \mathbb{R}^2 \rightarrow \mathbb{R}$$
$$x \mapsto \min_{p \in \mathcal{P}} \|x - p\|$$

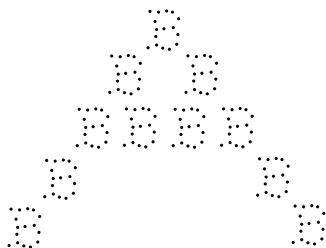


Dendrogram \rightarrow barcode

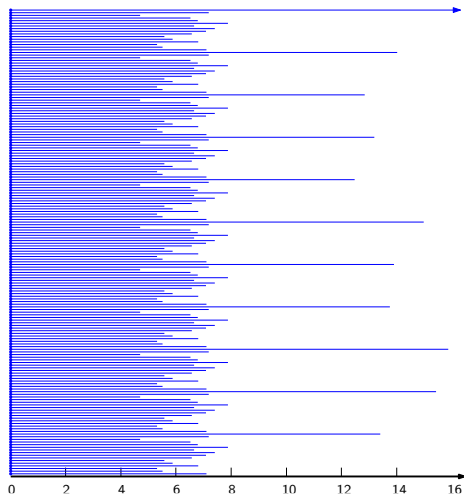


Dendrogram, Persistence Barcodes and Generalization

$$d_{\mathcal{P}} : \mathbb{R}^2 \rightarrow \mathbb{R}$$
$$x \mapsto \min_{p \in \mathcal{P}} \|x - p\|$$

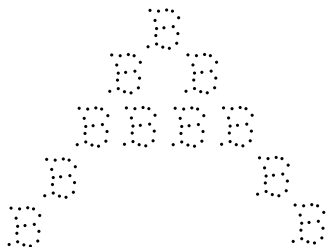


Dendrogram \rightarrow barcode



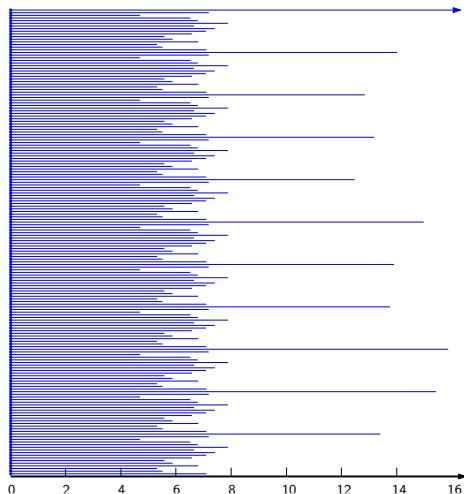
Dendrogram, Persistence Barcodes and Generalization

$$d_{\mathcal{P}} : \mathbb{R}^2 \rightarrow \mathbb{R}$$
$$x \mapsto \min_{p \in \mathcal{P}} \|x - p\|$$



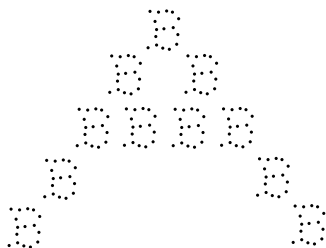
Barcode is:

- less informative
- more stable



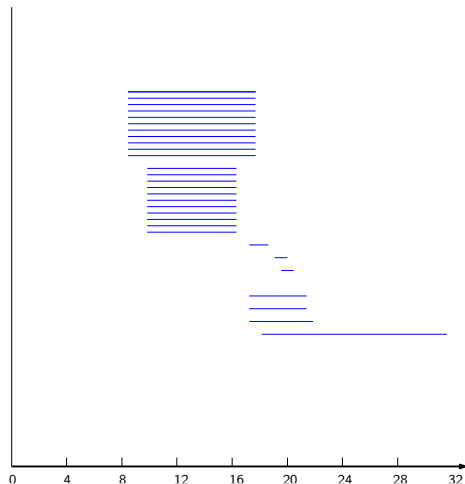
Dendrogram, Persistence Barcodes and Generalization

$$d_{\mathcal{P}} : \mathbb{R}^2 \rightarrow \mathbb{R}$$
$$x \mapsto \min_{p \in \mathcal{P}} \|x - p\|$$



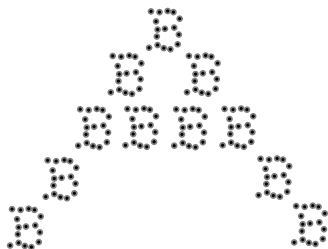
Barcode is:

- less informative
- more stable
- generalizable



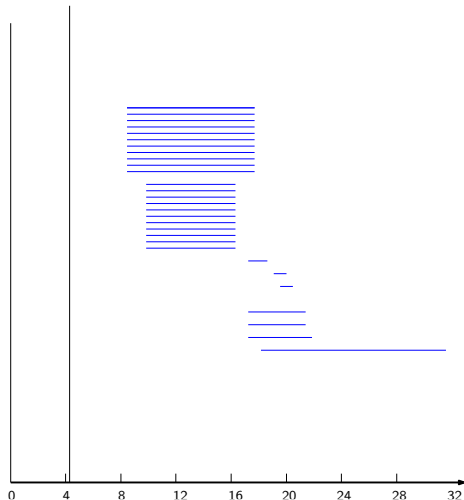
Dendrogram, Persistence Barcodes and Generalization

$$d_{\mathcal{P}} : \mathbb{R}^2 \rightarrow \mathbb{R}$$
$$x \mapsto \min_{p \in \mathcal{P}} \|x - p\|$$



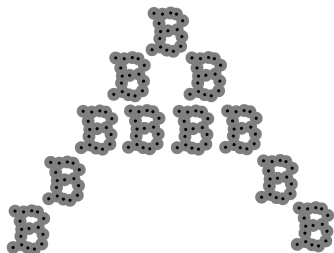
Barcode is:

- less informative
- more stable
- generalizable



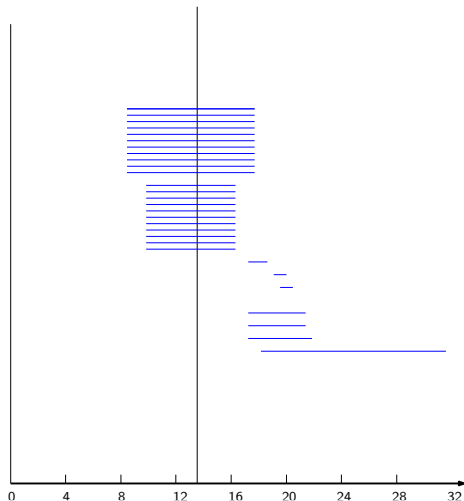
Dendrogram, Persistence Barcodes and Generalization

$$d_{\mathcal{P}} : \mathbb{R}^2 \rightarrow \mathbb{R}$$
$$x \mapsto \min_{p \in \mathcal{P}} \|x - p\|$$



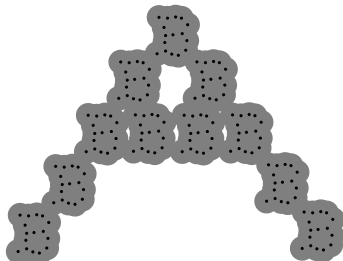
Barcode is:

- less informative
- more stable
- generalizable



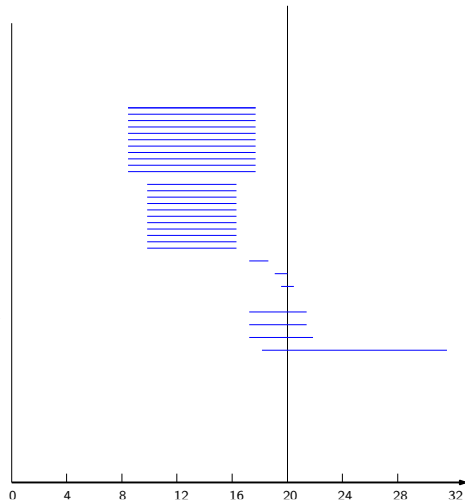
Dendrogram, Persistence Barcodes and Generalization

$$d_{\mathcal{P}} : \mathbb{R}^2 \rightarrow \mathbb{R}$$
$$x \mapsto \min_{p \in \mathcal{P}} \|x - p\|$$



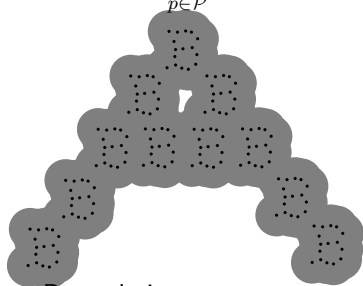
Barcode is:

- less informative
- more stable
- generalizable



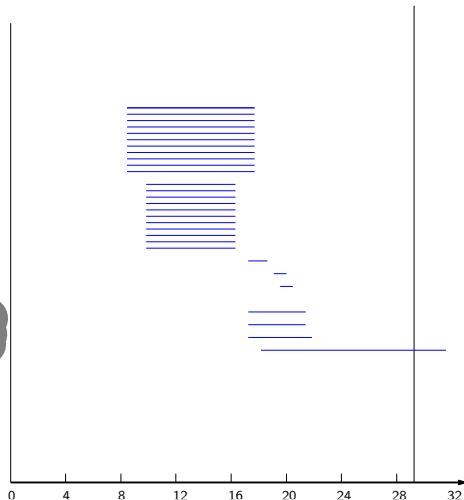
Dendrogram, Persistence Barcodes and Generalization

$$d_{\mathcal{P}} : \mathbb{R}^2 \rightarrow \mathbb{R}$$
$$x \mapsto \min_{p \in \mathcal{P}} \|x - p\|$$



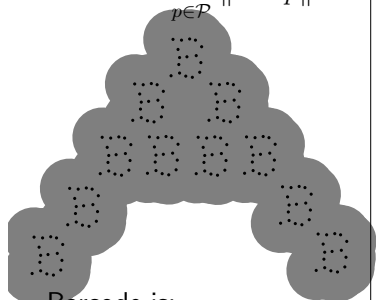
Barcode is:

- less informative
- more stable
- generalizable



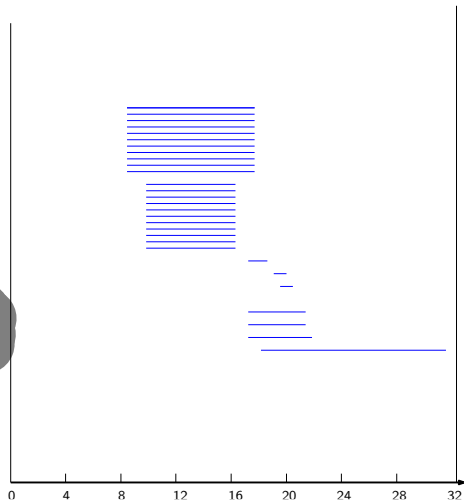
Dendrogram, Persistence Barcodes and Generalization

$$d_{\mathcal{P}} : \mathbb{R}^2 \rightarrow \mathbb{R}$$
$$x \mapsto \min_{p \in \mathcal{P}} \|x - p\|$$



Barcode is:

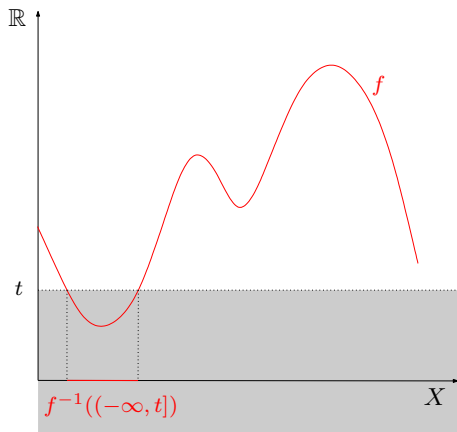
- less informative
- more stable
- generalizable



Persistence Diagrams

Inside the black box:

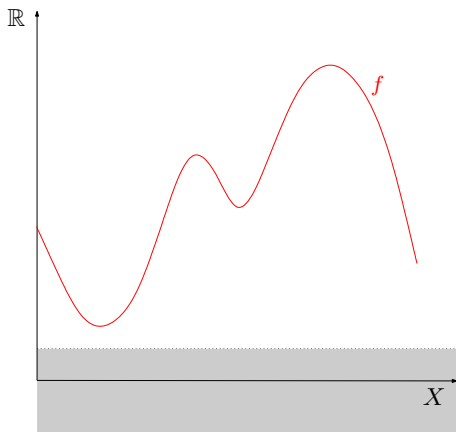
- Nested family (filtration) of sublevel-sets $f^{-1}((-\infty, t])$, for $t \in \mathbb{R}$.
- Track the evolution of the topology (homology) of the family.



Persistence Diagrams

Inside the black box:

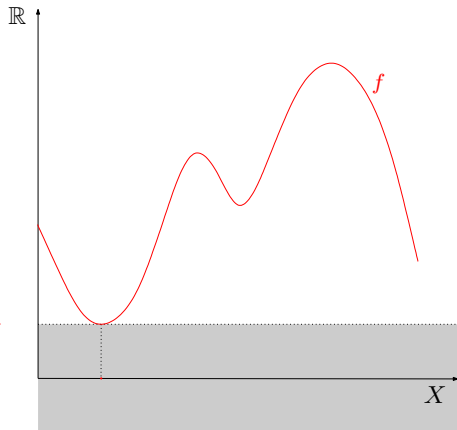
- Nested family (filtration) of sublevel-sets $f^{-1}((-\infty, t])$, for $t \in \mathbb{R}$.
- Track the evolution of the topology (homology) of the family.



Persistence Diagrams

Inside the black box:

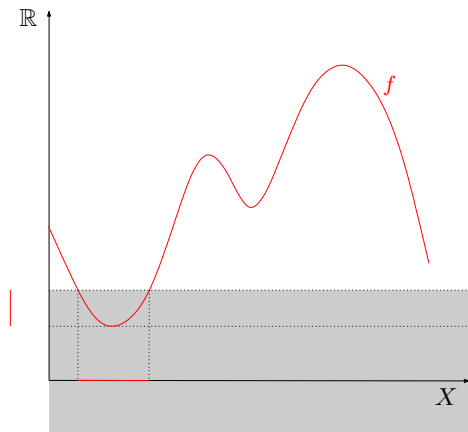
- Nested family (filtration) of sublevel-sets $f^{-1}((-\infty, t])$, for $t \in \mathbb{R}$.
- Track the evolution of the topology (homology) of the family.



Persistence Diagrams

Inside the black box:

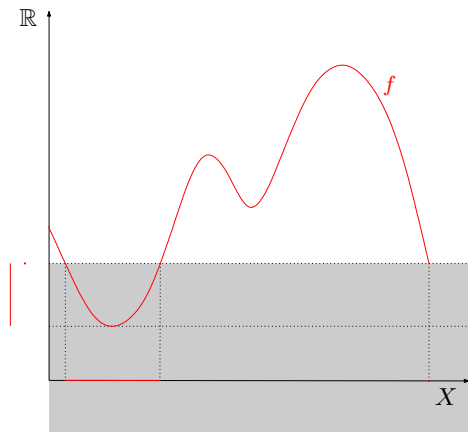
- Nested family (filtration) of sublevel-sets $f^{-1}((-\infty, t])$, for $t \in \mathbb{R}$.
- Track the evolution of the topology (homology) of the family.



Persistence Diagrams

Inside the black box:

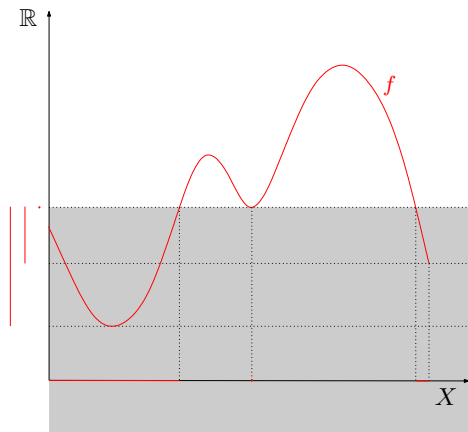
- Nested family (filtration) of sublevel-sets $f^{-1}((-\infty, t])$, for $t \in \mathbb{R}$.
- Track the evolution of the topology (homology) of the family.



Persistence Diagrams

Inside the black box:

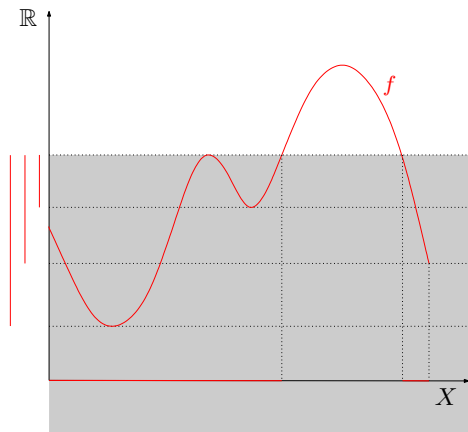
- Nested family (filtration) of sublevel-sets $f^{-1}((-\infty, t])$, for $t \in \mathbb{R}$.
- Track the evolution of the topology (homology) of the family.



Persistence Diagrams

Inside the black box:

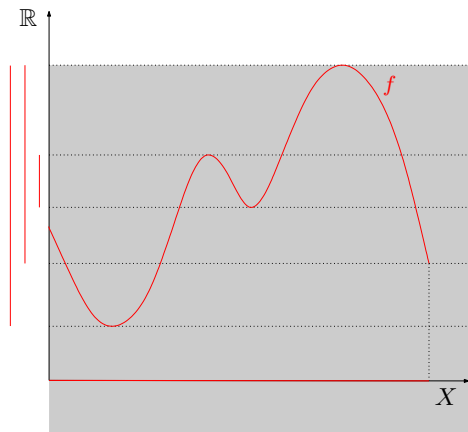
- Nested family (filtration) of sublevel-sets $f^{-1}((-\infty, t])$, for $t \in \mathbb{R}$.
- Track the evolution of the topology (homology) of the family.



Persistence Diagrams

Inside the black box:

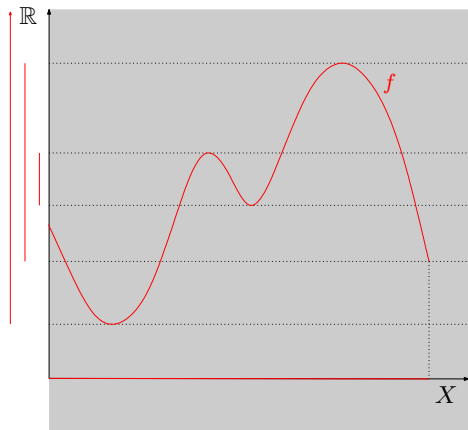
- Nested family (filtration) of sublevel-sets $f^{-1}((-\infty, t])$, for $t \in \mathbb{R}$.
- Track the evolution of the topology (homology) of the family.



Persistence Diagrams

Inside the black box:

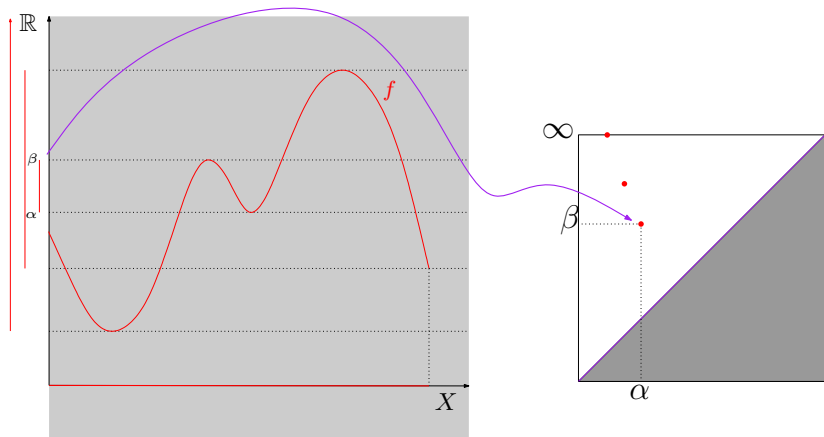
- Nested family (filtration) of sublevel-sets $f^{-1}((-\infty, t])$, for $t \in \mathbb{R}$.
- Track the evolution of the topology (homology) of the family.



Persistence Diagrams

Inside the black box:

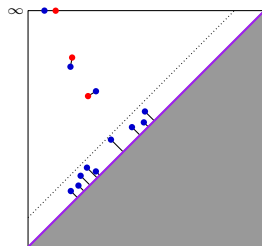
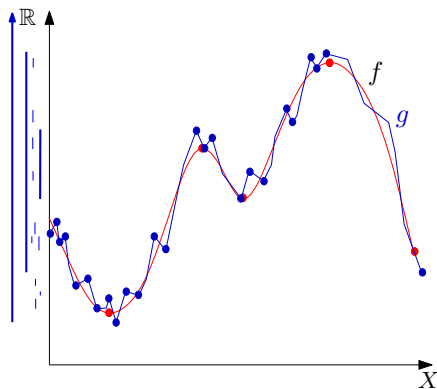
- Nested family (filtration) of sublevel-sets $f^{-1}((-\infty, t])$, for $t \in \mathbb{R}$.
- Track the evolution of the topology (homology) of the family.



Persistence Diagrams

Inside the black box:

- Nested family (filtration) of sublevel-sets $f^{-1}((-\infty, t])$, for $t \in \mathbb{R}$.
- Track the evolution of the topology (homology) of the family.

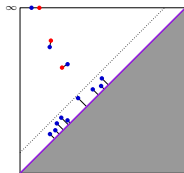
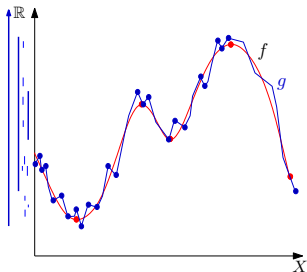


Persistence Diagrams

Definition (Bottleneck Distance)

Given two diagrams F and G ,

$d_b(F, G) = \inf\{\delta \mid \text{there exists a } \delta\text{-correspondence between } F \text{ and } G\}$.

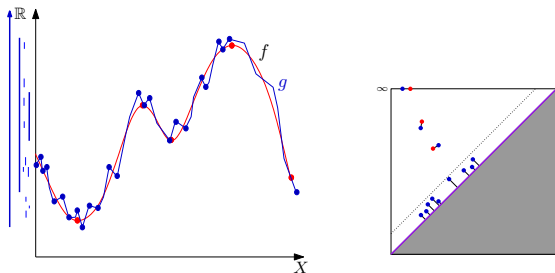


Persistence Diagrams

Definition (Bottleneck Distance)

Given two diagrams F and G ,

$d_b(F, G) = \inf\{\delta \mid \text{there exists a } \delta\text{-correspondence between } F \text{ and } G\}$.



Theorem (Stability of Persistence)

For all nice functions $f, g : X \rightarrow \mathbb{R}$,

$$d_b(\text{dgm}(f), \text{dgm}(g)) \leq \|f - g\|_\infty.$$

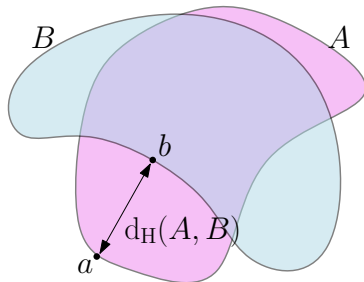
Stability for Sets

Definition (Hausdorff Distance)

The **Hausdorff distance** between two compact sets A and $B \subset \mathbb{R}^D$ is

$$d_H(A, B) = \|d_A(\cdot) - d_B(\cdot)\|_\infty,$$

where $d_K(x) = \inf_{p \in K} \|x - p\|$ is the distance to K .



Stability for Sets

Definition (Hausdorff Distance)

The **Hausdorff distance** between two compact sets A and $B \subset \mathbb{R}^D$ is

$$d_H(A, B) = \|d_A(\cdot) - d_B(\cdot)\|_\infty,$$

where $d_K(x) = \inf_{p \in K} \|x - p\|$ is the distance to K .

Proposition (Persistence Stability for Sets)

Write $\text{dgm}(K)$ for the diagram of the offset filtration

$$K^r = d_K^{-1}([0, r]), \text{ for } r \geq 0.$$

Then for all compact $A, B \subset \mathbb{R}^D$,

$$d_b(\text{dgm}(A), \text{dgm}(B)) \leq \|d_A(\cdot) - d_B(\cdot)\|_\infty = d_H(A, B).$$

Stability for Sets

Definition (Hausdorff Distance)

The **Hausdorff distance** between two compact sets A and $B \subset \mathbb{R}^D$ is

$$d_H(A, B) = \|d_A(\cdot) - d_B(\cdot)\|_\infty,$$

where $d_K(x) = \inf_{p \in K} \|x - p\|$ is the distance to K .

Proposition (Persistence Stability for Sets)

Write $\text{dgm}(K)$ for the diagram of the offset filtration

$$K^r = d_K^{-1}([0, r]), \text{ for } r \geq 0.$$

Then for all compact $A, B \subset \mathbb{R}^D$,

$$d_b(\text{dgm}(A), \text{dgm}(B)) \leq \|d_A(\cdot) - d_B(\cdot)\|_\infty = d_H(A, B).$$

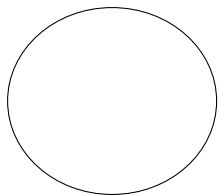
**Approximating persistence reduces to approximating sets
for Hausdorff loss.**

Homology in a Nutshell

β_0 : connected components

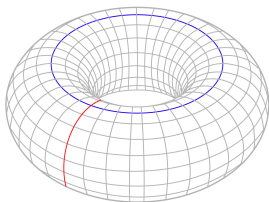
β_1 : holes

β_2 : voids



$$\beta_0 = 1$$

$$\beta_1 = 1$$



$$\beta_0 = 1$$

$$\beta_1 = 2$$

$$\beta_2 = 1$$



$$\beta_0 = 1$$

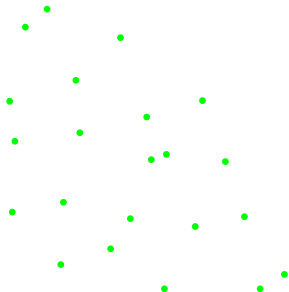
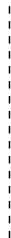
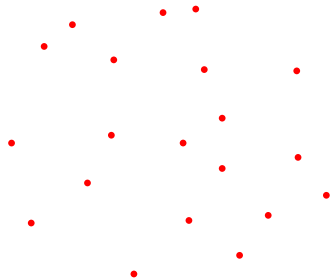
$$\beta_1 = 1$$

$$\beta_2 = 0$$

Support Estimation

Data: A n -sample $X_1, \dots, X_n \sim_{i.i.d.} P$.

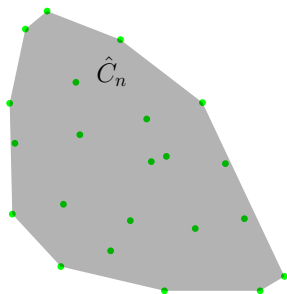
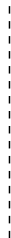
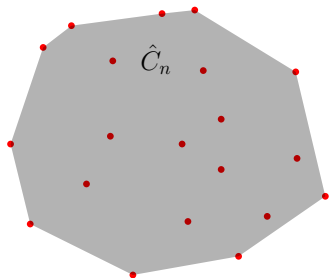
Goal: Estimate the set $C = \text{Support}(P) = \bigcap_{\substack{K \subset \mathbb{R}^D \text{ closed} \\ P(K)=1}} K$.



Support Estimation

Data: A n -sample $X_1, \dots, X_n \sim_{i.i.d.} P$.

Goal: Estimate the set $C = \text{Support}(P) = \bigcap_{\substack{K \subset \mathbb{R}^D \text{ closed} \\ P(K)=1}} K$.



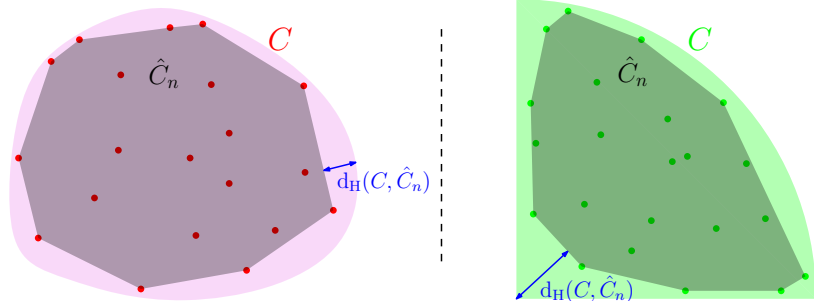
If we know (by advance) that C is convex, a good candidate is

$$\hat{C}_n = \text{Conv}(\{X_1, \dots, X_n\}).$$

Support Estimation

Data: A n -sample $X_1, \dots, X_n \sim_{i.i.d.} P$.

Goal: Estimate the set $C = \text{Support}(P) = \bigcap_{\substack{K \subset \mathbb{R}^D \text{ closed} \\ P(K)=1}} K$.



If we know (by advance) that C is convex, a good candidate is

$$\hat{C}_n = \text{Conv}(\{X_1, \dots, X_n\}).$$

Support Estimation: Convex Case(s)

Theorem (Dümbgen, Walther – 1996)

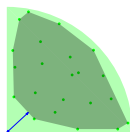
Assume that $P = \text{Unif}_C$ is uniform over the convex set $C \subset \mathbb{R}^D$.

Write

$$\hat{C}_n = \text{Conv}(\{X_1, \dots, X_n\}).$$

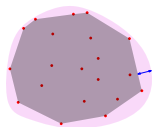
– Then,

$$d_H(C, \mathbb{X}_n) \leq d_H(C, \hat{C}_n) = O\left(\frac{\log n}{n}\right)^{\frac{1}{D}} \quad a.s.$$

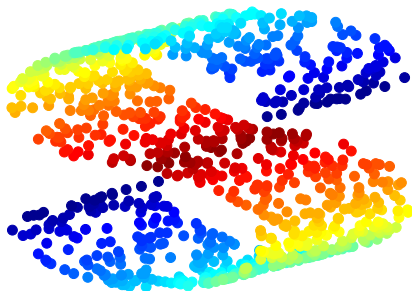


– If in addition, ∂C is \mathcal{C}^2 ,

$$d_H(C, \hat{C}_n) = O\left(\frac{\log n}{n}\right)^{\frac{2}{D+1}} \quad a.s.$$



Beyond Convexity



How to model the support of these data?

- Low-dimensional and curved \rightarrow Submanifold of \mathbb{R}^D .
- Not convex, but locally around it the projection uniquely defined.

Reminder: For a closed set $C \subset \mathbb{R}^D$,

$C \subset \mathbb{R}^D$ is convex \Leftrightarrow Every $z \in \mathbb{R}^D$ has a unique nearest neighbor on C
i.e. $\exists!$ $\pi_C(z) \in C$ with $\|z - \pi_C(z)\| = d_C(z)$.

Medial Axis

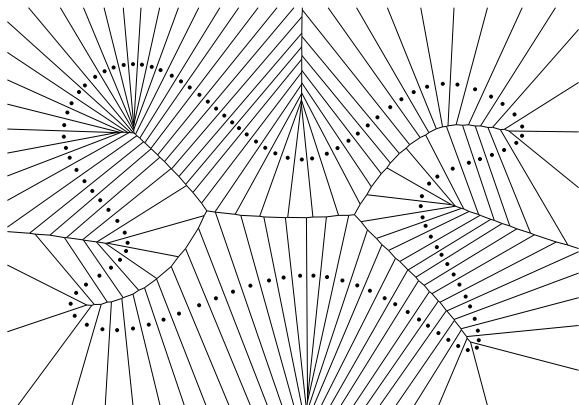
The **medial axis** of $M \subset \mathbb{R}^D$ is the set of points that have at least two nearest neighbors on M .

$$\text{Med}(M) = \{z \in \mathbb{R}^D, z \text{ has several nearest neighbors on } M\},$$

Medial Axis

The **medial axis** of $M \subset \mathbb{R}^D$ is the set of points that have at least two nearest neighbors on M .

$$\text{Med}(M) = \{z \in \mathbb{R}^D, z \text{ has several nearest neighbors on } M\},$$

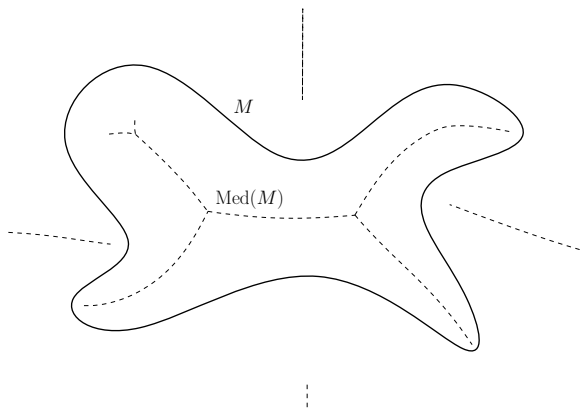


Medial axis of a point cloud (Voronoi faces)

Medial Axis

The **medial axis** of $M \subset \mathbb{R}^D$ is the set of points that have at least two nearest neighbors on M .

$$\text{Med}(M) = \{z \in \mathbb{R}^D, z \text{ has several nearest neighbors on } M\},$$



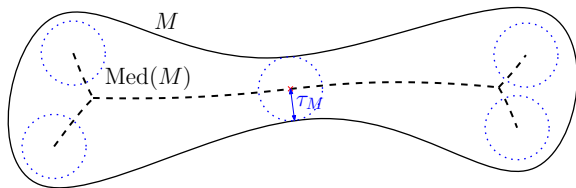
Medial axis of a continuous subset

Reach

For a closed subset $M \subset \mathbb{R}^D$, the **reach** τ_M of M is the least distance to its medial axis:

$$\tau_M = \inf_{x \in M} d_{\text{Med}(M)}(x),$$

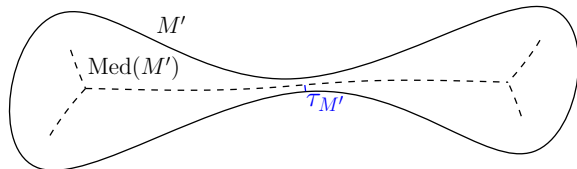
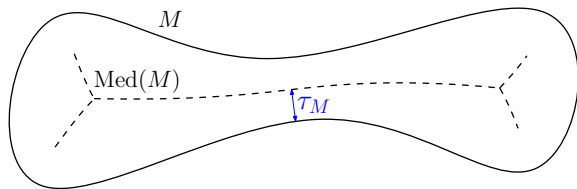
where for all $x \in \mathbb{R}^D$, $d_K(x) = \inf_{p \in K} \|x - p\|$.



One can also flip the formula:

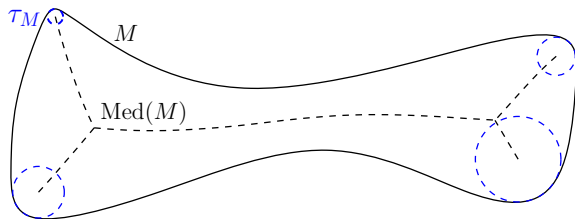
$$\tau_M = \inf_{z \in \text{Med}(M)} d_M(z).$$

Global Regularity



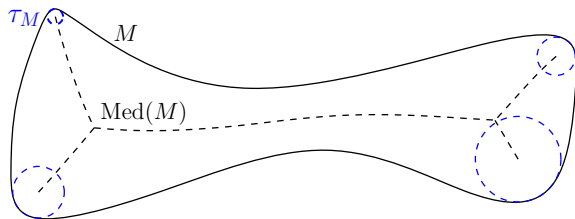
Narrow bottleneck structure $\Rightarrow \tau_M \ll 1$.

Local Regularity



High curvature \Leftrightarrow Small radius of curvature $\Rightarrow \tau_M \ll 1$.

Local Regularity



High curvature \Leftrightarrow Small radius of curvature $\Rightarrow \tau_M \ll 1$.

Proposition (Federer – 1959, Niyogi *et al.* – 2006)

Let II_x^M denote the second fundamental form of M .

For all unit tangent vector $v \in T_x M$, $\|II_x^M(v, v)\| \leq 1/\tau_M$.

As a consequence, the sectional curvatures κ of M satisfy

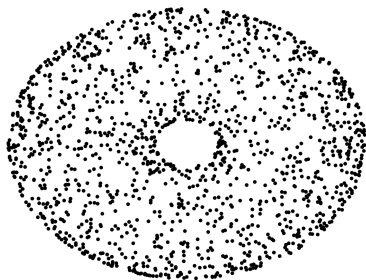
$$-2/\tau_M^2 \leq \kappa \leq 1/\tau_M^2.$$

Statistical Model

$X_1, \dots, X_n \stackrel{i.i.d.}{\sim} P$, where $M = \text{Support}(P) \subset \mathbb{R}^D$ satisfies:

- M is a compact connected d -dimensional submanifold,
- M has no boundary,
- $\tau_M \geq \tau_{\min} > 0$,
- P is (almost) the uniform distribution on M .

The set of distributions satisfying these conditions is denoted by \mathcal{P} .



A Reconstruction Theorem

Theorem (A, Levrard – 2018)

There exists a computable estimator \hat{M} such that for all $n \geq 1$,

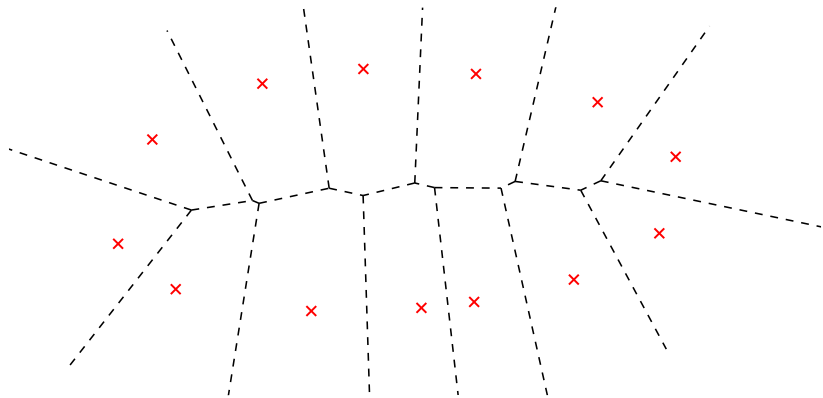
$$\mathbb{E}_{P^n} \left[d_H(M, \hat{M}) \right] \leq C \left(\frac{\log n}{n} \right)^{2/d},$$

where $C = C_{\tau_{\min}, d}$ does not depend on the ambient dimension D .

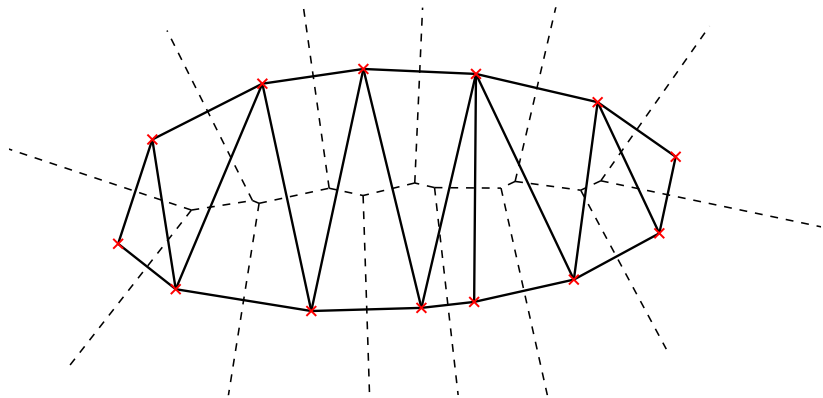
Outline of the Method



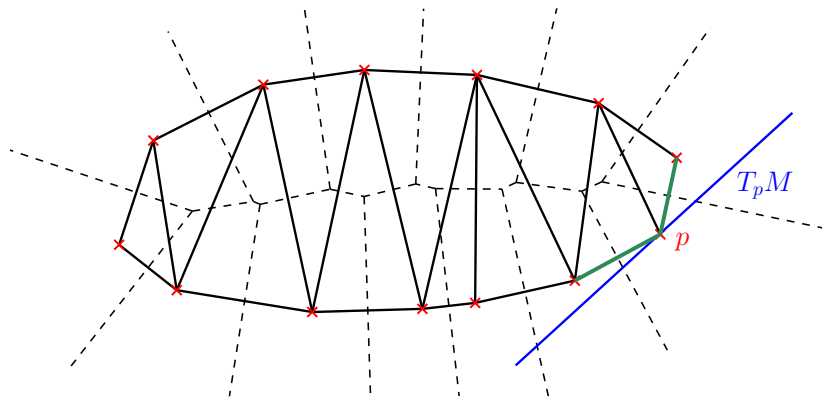
Outline of the Method



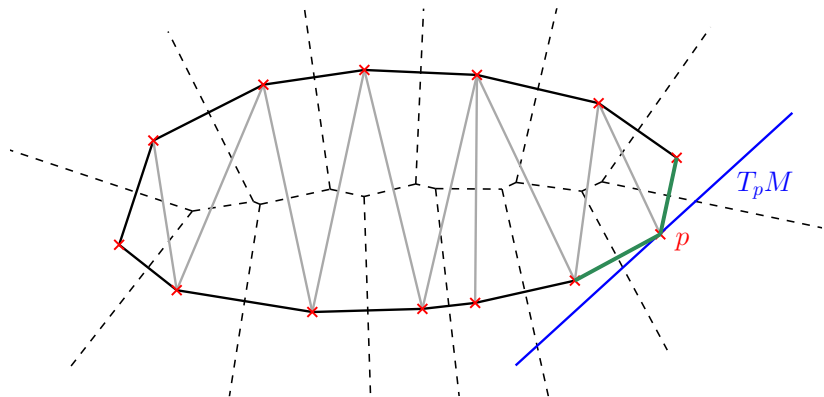
Outline of the Method



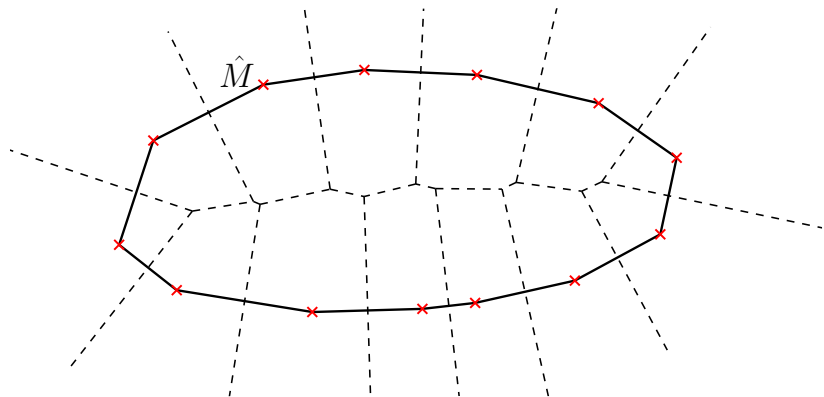
Outline of the Method



Outline of the Method



Outline of the Method



The Tangential Delaunay Complex [Boissonnat & Ghosh – 2014]

Optimality: Studying the Minimax Risk

The **minimax risk** over the statistical model \mathcal{P} is

$$\inf_{\hat{M}_n} \sup_{P \in \mathcal{P}} \mathbb{E}_{P^n} \left[d_H(M, \hat{M}_n) \right],$$

where the infimum is taken over all the estimators $\hat{M}_n = \hat{M}_n(\mathbb{X}_n)$ computed over a n -sample $\mathbb{X}_n = \{X_1, \dots, X_n\}$.

Optimality: Studying the Minimax Risk

The **minimax risk** over the statistical model \mathcal{P} is

$$\inf_{\hat{M}_n} \sup_{P \in \mathcal{P}} \mathbb{E}_{P^n} \left[d_H(M, \hat{M}_n) \right],$$

where the infimum is taken over all the estimators $\hat{M}_n = \hat{M}_n(\mathbb{X}_n)$ computed over a n -sample $\mathbb{X}_n = \{X_1, \dots, X_n\}$.

Proposition (Genovese *et al* – 2012)

For n large enough,

$$\inf_{\hat{M}_n} \sup_{P \in \mathcal{P}} \mathbb{E}_{P^n} \left[d_H(M, \hat{M}_n) \right] \leq C \left(\frac{\log n}{n} \right)^{\frac{2}{d}},$$

where $C = C_{d, \tau_{\min}}$

Optimality: Studying the Minimax Risk

The **minimax risk** over the statistical model \mathcal{P} is

$$\inf_{\hat{M}_n} \sup_{P \in \mathcal{P}} \mathbb{E}_{P^n} \left[d_H(M, \hat{M}_n) \right],$$

where the infimum is taken over all the estimators $\hat{M}_n = \hat{M}_n(\mathbb{X}_n)$ computed over a n -sample $\mathbb{X}_n = \{X_1, \dots, X_n\}$.

Proposition (Genovese *et al* – 2012)

For n large enough, (+ mild technical assumptions)

$$c \left(\frac{1}{n} \right)^{\frac{2}{d}} \leq \inf_{\hat{M}_n} \sup_{P \in \mathcal{P}} \mathbb{E}_{P^n} \left[d_H(M, \hat{M}_n) \right] \leq C \left(\frac{\log n}{n} \right)^{\frac{2}{d}},$$

where $C = C_{d, \tau_{\min}}$ and $c = c_{\tau_{\min}}$.

Lower Bound Technique: Le Cam's Lemma

Theorem (L. Le Cam)

For all $P_0, P_1 \in \mathcal{P}$,

$$\inf_{\hat{M}_n} \sup_{P \in \mathcal{P}} \mathbb{E}_{P^n} [d_H(M, \hat{M}_n)] \geq \frac{1}{2} d_H(M_0, M_1) (1 - \text{TV}(P_0, P_1))^n,$$

where

$$\text{TV}(P_0, P_1) = \sup_{B \in \mathcal{B}(\mathbb{R}^D)} |P_0(B) - P_1(B)|$$

denotes the total variation distance between P_0 and P_1 .

Lower Bound Technique: Le Cam's Lemma

Theorem (L. Le Cam)

For all $P_0, P_1 \in \mathcal{P}$,

$$\inf_{\hat{M}_n} \sup_{P \in \mathcal{P}} \mathbb{E}_{P^n} [\mathrm{d}_H(M, \hat{M}_n)] \geq \frac{1}{2} \mathrm{d}_H(M_0, M_1) (1 - \mathrm{TV}(P_0, P_1))^n,$$

where

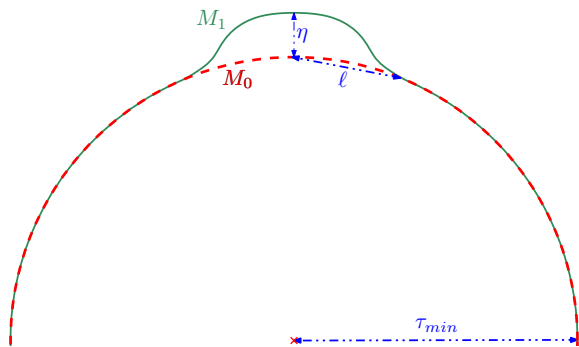
$$\mathrm{TV}(P_0, P_1) = \sup_{B \in \mathcal{B}(\mathbb{R}^D)} |P_0(B) - P_1(B)|$$

denotes the total variation distance between P_0 and P_1 .

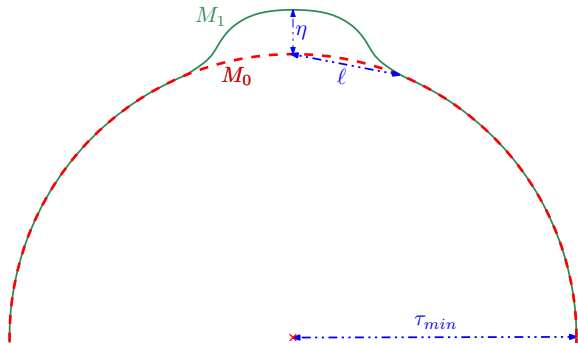
Deriving a good lower bound amounts to find P_0, P_1 such that:

- $P_0, P_1 \in \mathcal{P}$,
- $\mathrm{d}_H(M_0, M_1)$ is large,
- $\mathrm{TV}(P_0, P_1)$ is small.

Le Cam's Lemma Heuristic

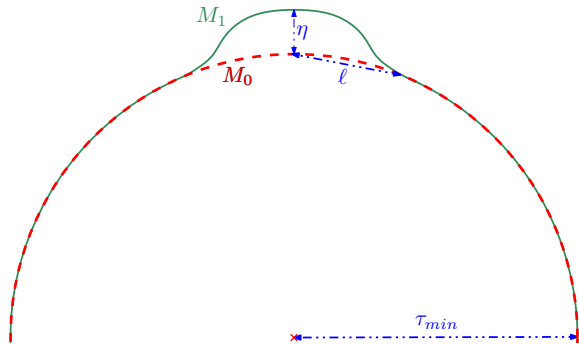


Le Cam's Lemma Heuristic



- P_0 and P_1 both belong to \mathcal{P} as soon as $\eta \lesssim \ell^2$,
- $d_H(M_0, M_1) \geq \eta$,
- $\text{TV}(P_0, P_1) \lesssim \ell^d$.

Le Cam's Lemma Heuristic

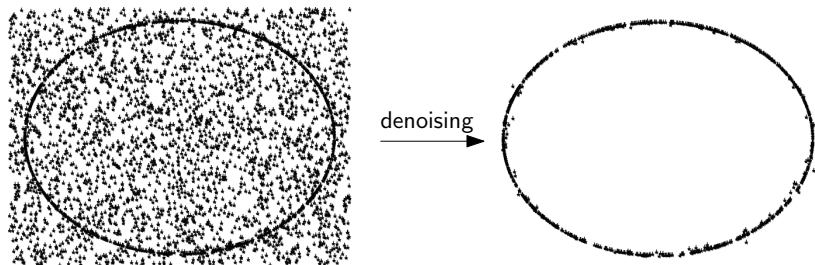


- P_0 and P_1 both belong to \mathcal{P} as soon as $\eta \lesssim \ell^2$,
- $d_H(M_0, M_1) \geq \eta$,
- $\text{TV}(P_0, P_1) \lesssim \ell^d$.

Hence, for $\eta \approx \ell^2$ and $\ell \approx (1/n)^{1/d}$,

$$\inf_{\hat{\tau}_n} \sup_{P \in \mathcal{P}} \mathbb{E}_{P^n} [d_H(M, \hat{M}_n)] \gtrsim \eta (1 - \ell^d)^n \approx (1/n)^{2/d}.$$

Extension to a Noisy Model

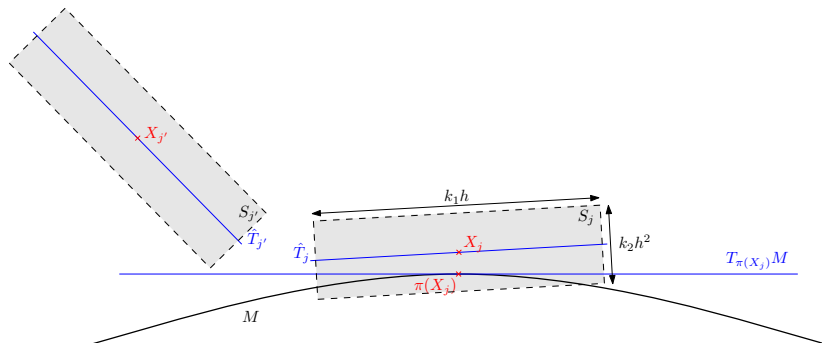


Theorem (A, Levrard – 2018)

For all $\delta > 0$, there exists a computable estimator $\hat{M}_n^{(\delta)}$ such that for all $n \geq 1$,

$$\mathbb{E}[\mathrm{d}_H(M, \hat{M}_n^{(\delta)})] \leq C \left(\frac{\log n}{n} \right)^{2/d-\delta}.$$

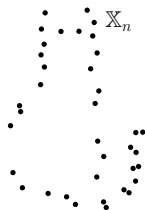
Denoising Outline



$$\begin{aligned}
 P(S(x, T_{\pi(x)}M)) &\asymp h^d && \text{if } d(x, M) \leq h^2, \\
 P(S(x, T)) &\asymp h^{2D-d} && \text{for all } T, \text{ if } d(x, M) > h^2,
 \end{aligned}$$

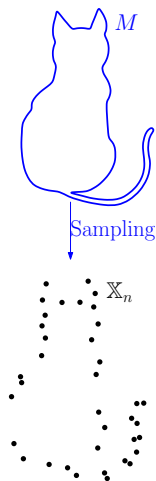
Since $h^{2D-d} \ll h^d$, the measure $P(S(x, T))$ of the slabs are discriminative for denoising.

The Catchy Slide...



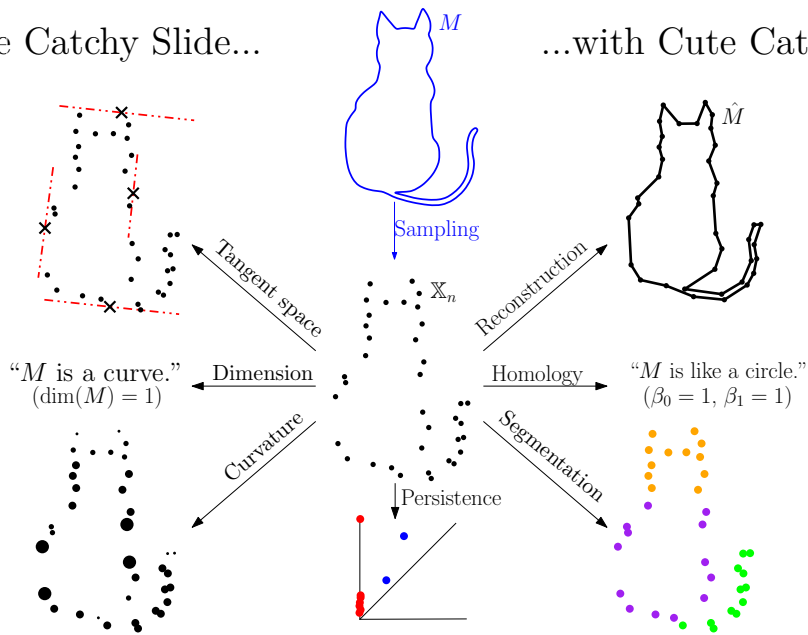
The Catchy Slide...

...with Cute Cats



The Catchy Slide...

...with Cute Cats



Lots of theoretical related topics:

- High-Dimensional statistics
- Nonparametric statistics
- Time series
- Computational geometry
- Geometry processing
- Abstract algebra

With applications in

- Material science
- Image analysis
- Physical chemistry
- Cosmology
- Network analysis
- ...